

7.1: Learning about the World through Surveys

Definition.

A **population** is a group of objects or people we wish to study.

- A **parameter** is a numerical value describing some aspect of the population (e.g. means and proportions)
- A **census** is a survey of *every member* of a population

A **sample** is a subset of the population of interest.

- A **statistic** is a numerical value describing some aspect of the sample
- Statistics are sometimes called **estimators**
- A **statistical inference** is the science of drawing conclusions about a population based on observing only a small subset of that population.

Example. In February 2014, the Pew Research Center surveyed 1428 cell phone users in the United States who were married or in a committed partnership. The survey found that 25% of cell phone owners felt that their spouse or partner was distracted by their cell phone when they were together.

Identify the population and the sample

Identify the parameter and the statistic

Sample statistics and population parameters are represented using different symbols (English for statistics, Greek for population):

	Statistics	Parameters
mean	\bar{x}	μ
standard deviation	s	σ
proportion	\hat{p}	p

Example. The City of Los Angeles provides an open data set of response times for emergency vehicles. Each row of the data set represents an emergency vehicle that has been sent to a particular emergency. A random sample of 1000 of these rows shows that the mean response time was 8.25 minutes. In addition, the proportion of vehicles that were ambulances was 0.328.

Using correct notation, identify the data given above.

What can we conclude about the overall population?

Definition.

A method is **biased** if it tends to produce the wrong value.

- **Sampling bias** results from a sample that is not representative of the population.
- **Measurement bias** results from questions that do not produce a true answer.

Example. For the following scenarios, identify any bias:

Online reviews (Amazon, Yelp, etc.)

Asking if people support a ‘fat tax’ on non-diet sugary soft drinks

Gallup poll calling landline phones

Using a poorly written question (e.g. double negative)

Definition.

A **simple random sample (SRS)** is where subjects from a population are drawn *at random* and *without replacement*. With an SRS, each member of the population has an equally likely chance of being selected.

Nonresponse bias results from people refusing to respond to the survey.

Example. Perform an SRS of 3 people from the list below:

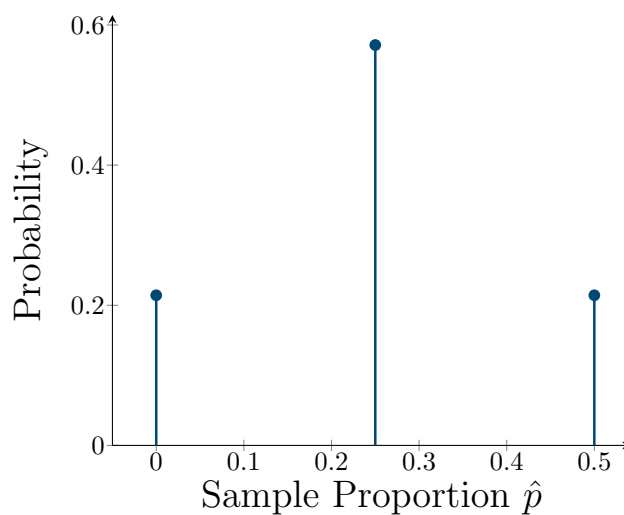
-
- 1 Alberto
 - 2 Justin
 - 3 Michael
 - 4 Audrey
 - 5 Brandy
 - 6 Nicole
-

7.2: Measuring the Quality of a Survey

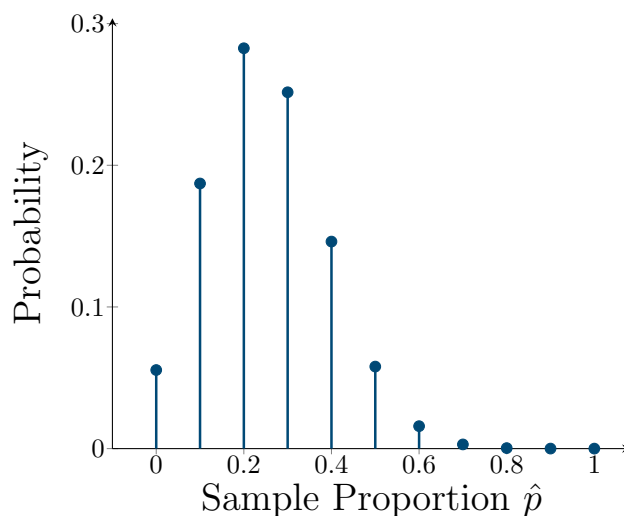
The true population proportion can be estimated by the sample proportion. How accurate can we expect our estimate to be?

- The *accuracy* of an estimation method is measured in terms of *bias*
- The *precision* of an estimation method is measured in terms of *standard error*

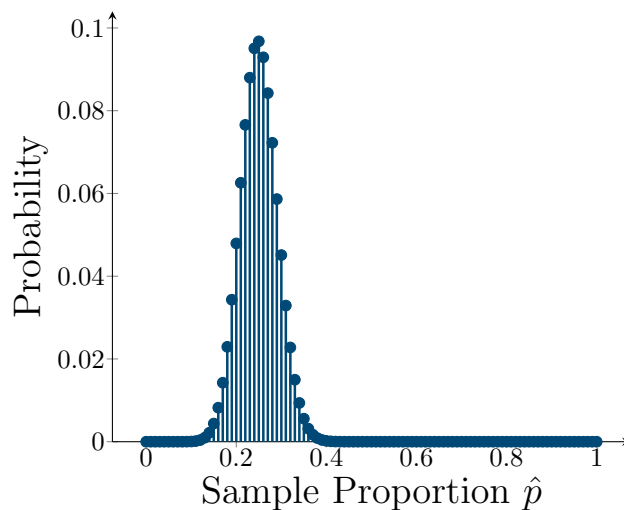
Example. Consider a group of 8 people, where 2 identify as female, and 6 identify as male. What is the true population proportion of females? When using a sample size of $n = 4$, what are possible sample proportions?



Example. Now, consider a group of 1000 people where 25% identify as female ($p = 0.25$). When using a sample size of $n = 10$, what are possible sample proportions?



Finally, consider a group of 1000 people where 25% identify as female ($p = 0.25$). When using a sample size of $n = 100$, what are possible sample proportions?



Definition.

- The **sampling distribution** is the probability distribution of \hat{p} .
- The **standard error** for \hat{p} is given by

$$SE = \sqrt{\frac{p(1-p)}{n}}$$

provided that

- The sample is randomly selected from the population of interest.
- If sampling without replacement, the population needs to be much larger than the sample size (e.g. at least 10 times bigger)

Since the true population proportion is typically unknown, we can estimate the standard error:

$$SE_{est} = \sqrt{\frac{\hat{p}(1-\hat{p})}{n}}$$

Note: Larger sample sizes have smaller standard error!

7.3: The Central Limit Theorem for Sample Proportions

Definition. (Central Limit Theorem (CLT))

When estimating a population proportion, p , if

1. *Random and Independent*: The sample is collected randomly from the population, and observations are independent of each other.
2. *Large Sample*: The sample size, n , is large enough that the sample can have at least 10 successes or failures.
3. *Big population*: If the sample is collected without replacement, then the population size must be at least 10 times bigger than the sample size.

then the sampling distribution for \hat{p} is approximately Normal, with mean p and standard deviation

$$SE = \sqrt{\frac{p(1-p)}{n}}.$$

This distribution is denoted as

$$N\left(p, \sqrt{\frac{p(1-p)}{n}}\right).$$

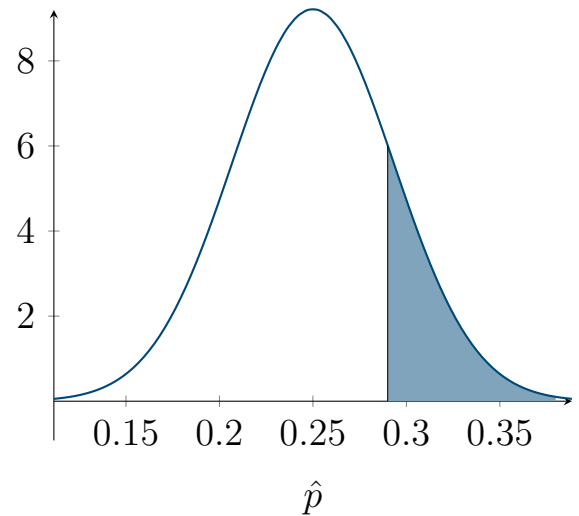
Example. Consider the groups from the previous section where $p = 0.25$ of the group identified as female. Suppose that $\hat{p} = 0.25$. If N represents the population size, and n the sample size, identify if the CLT can be applied.

$$N = 8, n = 4$$

$$N = 1000, n = 10$$

$$N = 1000, n = 100$$

Example. Consider the group of 1000 people where $p = 0.25$ identified as female. In a sample of $n = 100$ people, what is the probability that \hat{p} is at least 29%?



Example. Samuel Morse claimed that the true proportion of E's used in the English language is 0.12. Suppose we take a sample of 876 letters, and find a sample proportion of 0.1347. If we took another sample, what is the probability that the new sample proportion would be greater than 0.1347?

